

Honors Numerical Analysis

Jan 26, 2022

Course email / Google Group: hma22@nyu.edu

Today IEEE Arithmetic

Computers store information - namely numbers - in binary

format :

$$2 = (10)_2$$

$$7 = (111)_2 = (000111)_2$$

$$= 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0$$

$$= 0 \cdot 2^5 + \dots$$

Fractional numbers:

$$2.5 = (10.1)_2$$

$$2.25 = (10.01)_2$$

$$1 \cdot 2^1 + 0 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2}$$

Each 0 or 1 is a bit

Fixed representation : XXXX.XXXX ← Not how computers store info: unymore

Floating point representation : $\pm M \times 2^E$, $M \in [1, 2)$

(scientific/engineering notation)

↑ ↑
Mantissa Exponent

Example

$$23 = (10111)_2 \\ = (1.0111 \times 2^4)$$

Every multiplication by 2 moves decimal over to the right by 1 place.

$$\frac{1}{8} = (0.001)_2 \\ = (1.0 \times 2^{-3})$$

How many bits does a computer use?

"single precision" : 32 bits
1 bit for sign
8 bits for exponent
23 bits for mantissa

Ex: $5.5 = 1.011 \times 2^2$

$$= \left(\underbrace{0}_{\text{sign}} \mid \underbrace{E=00000010}_{\text{exponent}} \mid \underbrace{M=10110\dots0}_{\text{mantissa}} \right)$$

always 1 - could ignore, except for the number 0.
(hidden bit representation).

Def A floating point number is one that can be represented exactly using the above structure.

"double precision" : 64 bits
 1 bit for sign
 11 bits for exponent
 52 bits for mantissa 1.
52 bits

These representations were standardized in the 1980s by the IEEE. (Int. of Electrical and Electronics Engineers) (IEEE 754)

In Julia : $x = 7.0$
 $\text{typeof}(x) \rightarrow \text{Float64}$ (i.e. double precision)

$x = 7$
 $\text{typeof}(x) \rightarrow \text{Int64}$ (i.e. 64 bit = 8 byte integer).

To represent negative exponents E , the assumption is that E is offset by 1023

$$\Rightarrow x = \underbrace{\pm}_{1 \text{ bit}} \underbrace{M}_{52 \text{ bits}} \times 2^{E-1023}$$

\downarrow \searrow
 $1.0 \dots 0$ $1010 \dots 0$
52 bits 11 bits

In Julia:
 $\text{bitstring}(x)$ returns IEEE representation
 $\text{Inf}, \infty, \neq 0$

$$E = (0111111111)_2 \Rightarrow 2^0$$

$$E = (1000000000)_2 \Rightarrow 2^{1024-1023} = 2^1$$

The point of IEEE 754 is to have consistent rules so that code is portable!

Special values:

$$\text{Inf} = \frac{\infty}{0} \Rightarrow M=0, E=1\dots 1$$

$$\pm 0 \Rightarrow M=0, E=0\dots 0$$

$$\text{NaN} = \sqrt{-1} \Rightarrow M \neq 0, E=1\dots 1$$

Rounding

Example: $\frac{1}{10} = (0.1)_{10}$

$$= (1.1001100)_2 \times 2^{-4}$$

this must be rounded and stored as a floating point number.

There are 4 options: up, down, to-zero, (nearest)
 ↓ default

Details are not important, just know that it's a standard specified by IEEE.

The important idea is the error made in rounding:

$$\text{Absolute rounding error} = |\text{round}(x) - x|$$

$$\text{Relative rounding error} = \frac{|\text{round}(x) - x|}{|x|} \quad \left. \vphantom{\frac{|\text{round}(x) - x|}{|x|}} \right\} \text{similar to percentage error}$$

Most important rule/consequence of the standard:

IEEE says that it must be the case that:

$$\begin{aligned} \underbrace{\text{round}(a+b)}_{\text{exact addition}} &= \text{round}(a) \oplus \text{round}(b) \\ &\quad \uparrow \\ &\quad \text{computer addition} \\ &= (a+b)(1+\delta) \end{aligned}$$

where $|\delta| \leq \epsilon$, with $\epsilon = \underline{\text{machine precision}}$

Machine precision is the distance between 1 and the next floating point number:

$$\epsilon = \underbrace{(0|01\dots1|0\dots00)}_2 - \underbrace{(0|01\dots1|0\dots01)}_2 = \underbrace{(1.0\dots01)}_2 = 1 + 2^{-52}$$

$\Rightarrow \epsilon = 2^{-52}$ on double precision

How to calculate this? \Rightarrow Julia demo

The IEEE standards enforce rules on the relative accuracy

of $+x - \div$:

$$\text{round}(a+b) = (a+b)(1+\delta)$$

$$\Rightarrow \frac{|\text{round}(a+b) - (a+b)|}{|a+b|} \leq \epsilon$$

In terms of absolute accuracy:

$$|\text{round}(a+b) - (a+b)| = |a+b| |\delta|$$

$$\leq (|a| + |b|) \epsilon$$

$$\leq 2 \max(|a|, |b|) \epsilon$$

Julia comments:

- problems installing?

julia-lang.org

- make plot of something

- formatted output